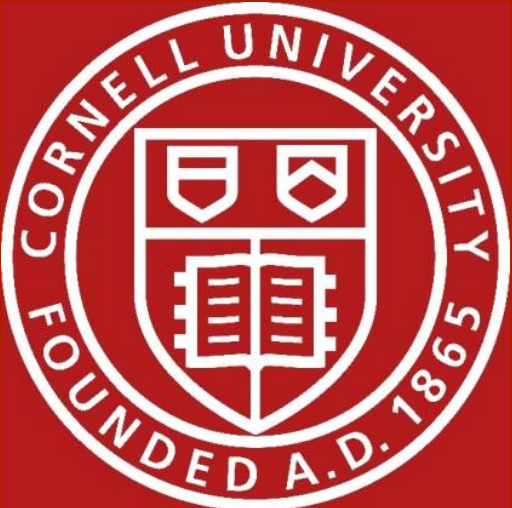


Incorporating W3C's DQV and PROV in CISER's Data Quality Review and Reproduction of Results Service



William C. Block Florio O. Arguillas Jeremy Williams
Cornell Institute for Social and Economic Research (CISER)

Corresponding author e-mail:
block@cornell.edu

Abstract

A year and a half after its implementation, CISER's Research Data Quality Review and Reproduction of Results Service (or R^2 , for short), a service developed to encourage sharing of high quality data, code, documentation, and metadata associated with a study for the purpose of reproducible research, continues to evolve and improve. This poster discusses: a) the service at its current state; b) the improvements made to the service including cost-reduction and buy-in strategies to encourage researchers to use the service; c) pre- and post-reproduction services to improve data, code, documentation, and metadata quality; d) the inclusion of the Cornell's CED²AR software for assessing and generating complete DDI metadata; and e) the utilization of the CISER Data Archive as the free and permanent home for the study and its associated files. Last, our poster will include for the first time our efforts to incorporate the W3C's Data Quality Vocabulary (DQV) and PROV metadata into our R^2 service. DQV provides a means to describe the subjective measures applied to ensure the integrity of data sources and introduces a way of expressing guidelines for data quality that producers of data can abide by. The W3C PROV ontology defines entities, agents, and activities related to the origin of resources, which is particularly important in the common scenario of referencing multiple datasets for a given investigation." By incorporating DQV and PROV into the R^2 data management workflow, we hope to provide a foundation of support improving the quality of research inputs and outputs, which in turn can have a chain effect for derivative data products.

		1996	
1988	Wanted	Mistimed	Unwanted
Wanted	CONSISTENT	INCONSISTENT	INCONSISTENT
Mistimed	INCONSISTENT	CONSISTENT	INCONSISTENT
Unwanted	INCONSISTENT	INCONSISTENT	CONSISTENT

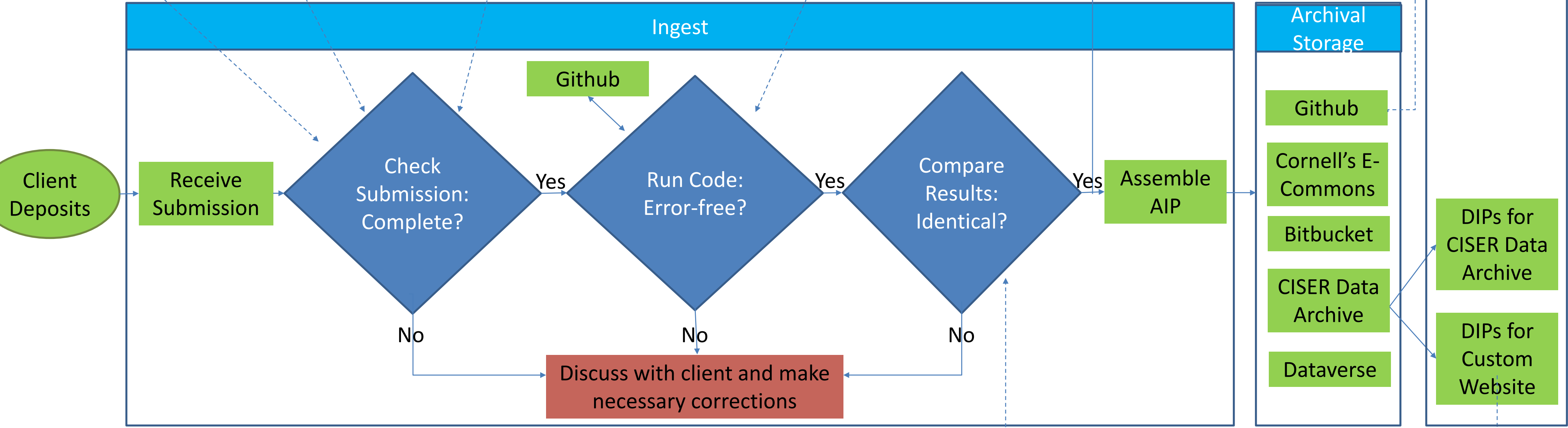
FIG. 1.—First Dependent Variable Construction

TABLE 1a			
CONSISTENCY OF WOMEN'S WANTEDNESS REPORTS BETWEEN TIME 1 AND TIME 2 ABOUT THE SAME PREGNANCY, ACCORDING TO PREGNANCY OUTCOME (WEIGHTED PERCENTAGES)			
Pregnancy Outcome	Live Births	Abortion	Stillbirth, miscarriage, ectopic pregnancy
Live birth	24.9	75.1	1338
Abortion	28.3	71.7	140
Stillbirth, miscarriage, ectopic pregnancy	32.9	67.0	240
Total	26.5	73.5	1718

TABLE 1b			
CONSISTENCY OF WOMEN'S PREGNANCY REPORTS BETWEEN TIME 1 AND TIME 2 ABOUT PARTNERS' VIEWS, ACCORDING TO PREGNANCY OUTCOME (WEIGHTED PERCENTAGES)			
Pregnancy Outcome	Live Births	Abortion	Stillbirth, miscarriage, ectopic pregnancy
Live birth	26.0	74.0	1260
Abortion	36.0	64.0	124
Stillbirth, miscarriage, ectopic pregnancy	33.1	66.9	218
Total	28.0	72.0	1602

ately for different Time 1 intention statuses. As a result, in those two tables the dependent variable can capture shifts toward more positive and more negative views. The dependent variable comprises three categories: More positive, more negative, and consistent. A more positive report resulted (1) if the woman claimed

Common problems: No variable and value labels



Cost-reduction strategies

- Data curation and management training
- Code writing and organization training
- Code efficiency training e.g., macro programming, SQL programming
- Version control software training e.g., Github



International Association for Social Science Information Services & Technology

Common problems:

- Some results do not match article
- Not all figures printed on the article are produced by the code. Some involved other software packages such as Excel.
- Order of variables in the model in the printed table do not match the order of the variables in the output table for that model, which slows down the verification process.

Golden Rule of CISER's Replication Service:

Output produced by running code against the data should be identical to the publication up to the last decimal place. Slight deviation is not acceptable and must be investigated.

Common problems:

- Very long, complex codes
- Unnecessary/excess sections of codes whose outputs are not found in the paper (this delays replication because the Staff has to go through the entire code and its output, and figure out where they are on the paper)
- Code points to subdirectories for retrieving or saving data, thus Staff has to recreate the directory structure for the code to run correctly
- Some codes are not efficient, but will not be modified by the Staff. The Staff, however, will suggest ways to make it efficient.
- Codes are often multiple files with no indication of sequence. Reproducer has to determine which to run first especially if codes build on top of the other.

